

Classification with Boosting of Extreme Learning Machine Over Arbitrarily Partitioned Data

Ferhat Özgür Çatak

04.21.2015 / Accepted: date

Abstract Machine learning based computational intelligence methods are widely used to analyze large scale data sets in this age of big data. Extracting useful predictive modeling from these types of data sets is a challenging problem due to their high complexity. Analyzing large amount of streaming data that can be leveraged to derive business value is another complex problem to solve. With high levels of data availability (*i.e. Big Data*) automatic classification of them has become an important and complex task. Hence, we explore the power of applying MapReduce based Distributed AdaBoosting of Extreme Learning Machine (ELM) to build a predictive bag of classification models. Accordingly, (i) data set ensembles are created; (ii) ELM algorithm is used to build weak learners (classifier functions); and (iii) builds a strong learner from a set of weak learners. We applied this training model to the benchmark knowledge discovery and data mining data sets.

Keywords Extreme Learning Machine · AdaBoost · Ensemble Methods · MapReduce

1 Introduction

It is clear that there has been an unexpected increase in the quantity and variety of data generated worldwide by computers, mobile phones, and sensors. Just as computer technology evolved, the quantity and variety of data has also increased, becoming more focused on storing every type of data, the so-called Big Data. As the volume of data to build a predictive model increases,

Ferhat Özgür Çatak
TÜBİTAK BİLGEM, Cyber Security Institute
Kocaeli, Turkey
Tel.: +90-262-6481000
Fax: +90-262-6481100
E-mail: ozgur.catak@tubitak.gov.tr

the complexity of that training increases too. As a result, building actionable predictive modeling of a large scale unstructured data set is a definitive Big Data problem. Predictive learning models try to discover patterns of training data and label new data instances to the correct output value. To efficiently handle unstructured large scale big data sets, it is critical to develop new machine learning methods that combine several boosting and classification algorithms.

Extreme Learning Machine (ELM) was proposed by [1] based on generalized Single-hidden Layer Feedforward Networks (SLFNs). Main characteristics of ELM are small training time compared to traditional gradient-based learning methods, high generalization property of predicting unseen examples with multi-class labels and parameter free with randomly generated hidden nodes. ELM algorithm is used in many different areas including document classification [2], bioinformatics [3] multimedia recognition [4, 5].

In recent years, much computational intelligence research has been devoted to building predictive modeling of distributed and parallel frameworks. In this research, the proposed learning model creates data chunks with varying size and bag of classifier functions using ELM algorithm trained with these arbitrary chosen sub data set with AdaBoosting method for large scale predictions. By creating data chunks from the training data set using the MapReduce paradigm, each subset of the training data set is used to find out the set of ELM ensembles as a single global classifier function.

The main objective of this work is to train large scale data sets using ELM and AdaBoost. Another objective is to achieve the model's classification performance with same or close to the conventional ELM method. Conventional ELM training cannot be applied to large scale data sets on a single computer because of their complexity. Then experiments section is split into two sub-sections: "commonly used data sets" in Section 5.1.1 and "large scale data sets" in Section 5.1.2. Commonly used data sets are suitable for training on a single computer with the conventional ELM algorithm. We trained these data sets both conventional and proposed methods to show the classification performance changes of the proposed method. Classification performance results are shown in Section 5.3.

The contributions of this paper are as follows:

- A generative MapReduce technique based AdaBoosted ELM classification model is proposed for learning, and thus, faster classification model training is achieved.
- This research proposes a new learning method for AdaBoosted ELM that achieves parallelization both in large scale data sets and reduced computational time of learning algorithm.
- Training computations of working nodes are independent from each other thus minimizing the data communication. The other approaches, including Support Vector Machine training need data communication for the support vector exchange. [6, 7, 8].

The rest of the paper is organized as follows: Section 2 briefly introduces some of the earlier works related to our problem. Section 3 describes algorithm ELM, AdaBoost and MapReduce technique. Section 4 and Section 5 evaluates the proposed learning model. Section 6 concludes this paper.

2 Related work

In this section, we describe the general overview of literature review. Section 2.1 describes the general distributed ELM methods. Section 2.2 shows the MapReduce based ELM training methods.

2.1 Literature Review Overview

MapReduce based learning algorithms from distributed data chunks has been studied by many researchers. Many different MapReduce based learning solutions over arbitrary partitioned data have been proposed recently. Some popular MapReduce based solutions to train machine learning algorithms in the literature include the following. Panda et al. proposed a learning tree model which is based on series of distributed computations, and implements each one using the MapReduce model of distributed computation [9]. Zhang et al. develops some algorithms using MapReduce to perform parallel data joins on large scale data sets [10]. Sun et al. use batch updating based hierarchical clustering to reduce computational time and data communication [11]. Their approach uses co-occurrence based feature selection to remove noisy features and decrease the dimension of the feature vectors. He et al. proposed parallel density based clustering algorithm (DBSCAN). They developed a partitioning strategy for large scale non-indexed data with a 4-stages MapReduce paradigm [12]. Zhao et al. proposed parallel k-means clustering based on MapReduce [13]. Their approaches focus on implementing k-means with the read-only convergence heuristic in the MapReduce pattern.

2.2 MapReduce Based ELM Training Methods

Section 2.2.1 - Section 2.2.5 describe five different MapReduce training methods of ELM algorithm.

2.2.1 *ELM**

Xin et al. proposed MapReduce based ELM training method called as ELM* [14]. Main idea behind this method is to calculate matrix multiplication of ELM to find weight vector. They show that Moore-Penrose generalized inverse operator is the most expensive computation part of the algorithm. As we know, matrix multiplication can be divide into smaller part. Using this property, they

proposed an efficient implementation of training phase to manage massive data sets. The final output of this method is a single classifier function. In this paper, they proposed two different versions of ELM*, naive and improved. In naive-ELM*, the algorithm has two classes, Class Mapper and Class Reducer. Both classes contain only one method. In improved ELM*, they decompose the calculation of matrix multiplication using MapReduce framework. Moreover, the proposed algorithm decreases the computation and communication cost. In the experimental platform, they used their synthetic data sets to evaluate the performance of the proposed algorithms with MapReduce framework.

2.2.2 OS-ELM based Classification in Hierarchical P2P Network

Sun et al. proposed OS-ELM [15] based distributed ensemble classification in P2P networks [16]. They apply the incremental learning principle of OS-ELM to hierarchical P2P network. They proposed two different versions of the ensemble classifier in hierarchical P2P, *one-by-one* ensemble classification and *parallel* ensemble classification. In *one-by-one* learning method, each peer, one by one, calculates the classifier with all the data. Therefore, this approach has a large network delay. In the *parallel* ensemble learning, all the classifiers are learnt from all the data in parallel manner. Conversely to ELM*, their experimental results are based on three different real data sets downloaded from the UCI repository.

2.2.3 Parallel online sequential ELM: POS-ELM

Wang et al. have been proposed parallel online sequential extreme learning machine (POS-ELM) method [17]. Main idea behind in this approach is to analyze the dependency relationships and the matrix calculations of OS-ELM [15]. Their experimental results are based on nine different real data sets downloaded from the UCI repository.

2.2.4 Distributed and Kernelized ELM: DK-ELM

Bi et al. have been proposed both distributed and kernelized ELM (DK-ELM) based on MapReduce [18]. The difference between ELM and Kernelized ELM is that K-ELM applies kernels opposite to create random feature mappings. They provide a distributed implementation RBF kernel matrix calculation in massive data learning applications. Their experimental results are based on four different real data sets downloaded from the UCI repository and four synthetic data sets.

2.2.5 ELM-MapReduce

Chen et al. have been proposed MapReduce based ELM ensemble classifier called ELM-MapReduce, for large scale land cover classification of remote sensing data [19]. Their approach contains two sequential phases: parallel training

of multiple ELM classifiers and voting mechanism. In parallel training phase of proposed method, each *Map* function computes an ELM classifier with a given training data set. In second phase called voting mechanism, a new MapReduce job is executed with a new partitioned test set into each *Map* function with notation $data_j$. In *Reduce* function of this phase, each $data_j$ is predicted with each ELM classifier trained in parallel training phase. Final classification predictions are the output of final *Reduce* function. Therefore, this approach has a high communication cost. Their experimental results are based synthetic remote sensing image of training data.

2.3 The Differences Between Proposed Model and Literature Review

The main differences are:

- In ELM \star , they use matrix multiplication decomposition. Each *Map* function is responsible to calculate the Moore-Penrose generalized inverse operation. And their method produces one single classifier. In the proposed model in our paper, each *Reduce* function produces ensemble classifier based on AdaBoost method. The final output ensemble classifier is a voting based combination of ensemble classifier trained in each *Reduce* phase.
- In OS-ELM based classification in hierarchical P2P Network, POS-ELM and DK-ELM, they propose ensemble classifier that combines multiple classifier trained with data chunks. Each peer classifier is learned from the local data. Therefore, each peer produces a single ELM classifier. In our method, each node (or peer) produces ensemble classifier to increase the classification accuracy.
- In ELM-MapReduce, they propose ensemble classifier with two different MapReduce jobs. In first MapReduce job, their approach produces a single ELM classifier in each *Map* function. In second MapReduce job, the test set is partitioned into each *Map* function and produces final predicted labels based on the voting mechanism of ELM classifiers that are trained in the first MapReduce job. In our method, prediction is not included, our aim is to create a final ensemble classifier in only one MapReduce job.

Table 1 shows the main differences of all proposed methods. There are five different columns that are *ensemble methods*, *single pass MapReduce*, *matrix multiplication*, *entire data set* and *network communication*. *Ensemble* column shows that the method builds a set of classifier function (i.e. ensemble model) to improve the accuracy performance of the final classification model. If an ensemble method is applied, then the performance of final model will have better accuracy result [20]. *Single Pass MapReduce* column shows that an iterative approach is not applied to the model. Entire learning phase is performed in a single pass of data through the job. *Matrix Multiplication* column shows the *hidden layer matrix* is calculated in each *Map* function. The hidden layer matrix computation is a compute intensive operation. *Entire Data Set* column shows each *Map* operation needs entire data set to build a final classifier model. *Network Communication* column shows that each *MapReduce* job

needs to communicate with another job. Network communication will affect negatively on training time of the algorithm.

Table 1: The Differences Between Proposed Model and Literature Review.

Method	Ensemble	Single Pass MapRe- duce	Matrix Multiplica- tion	Entire Data Set	Network Communi- cation
ELM*	No	Yes	No	Yes	No
OS-ELM	Yes	Yes	No	No	Yes
POS-ELM	Yes	Yes	No	Yes	No
DK-ELM	Yes	Yes	No	Yes	No
ELM-MapReduce	Yes	No	No	Yes	Yes
Proposed Method	Yes	Yes	No	No	No

3 Preliminaries

In this section, we introduce preliminary knowledge of ELM, AdaBoost and MapReduce briefly.

3.1 Extreme learning machine

ELM was originally proposed for the single-hidden layer feedforward neural networks [21, 22, 1]. Then, ELM was extended to the generalized single-hidden layer feedforward networks where the hidden layer may not be neuron like [23, 24]. The main advantages of the ELM classification algorithm are that ELM can be trained hundred times faster than traditional neural network or support vector machine algorithm since its input weights and hidden node biases are randomly created and output layer weights can be analytically calculated by using a least-squares method [25, 26]. The most noticeable feature of ELM is that its hidden layer parameters are selected randomly.

Given a set of training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}$, $\mathbf{x}_i \in \mathbb{R}^p$, $y_i \in \{1, 2, \dots, K\}$ sampled independently and identically distributed (i.i.d.) from some unknown distribution. The goal of a neural network is to learn a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ where \mathcal{X} is instance and \mathcal{Y} is the set of all possible labels. The output label of an single hidden-layer feedforward neural networks (SLFNs) with N hidden nodes can be described as

$$f_N(\mathbf{x}) = \sum_{i=1}^N \beta_i G(\mathbf{a}_i, b_i, \mathbf{x}), \mathbf{x} \in \mathbb{R}^n, \mathbf{a}_i \in \mathbb{R}^n \quad (1)$$

where \mathbf{a}_i and b_i are the learning parameters of hidden nodes and β_i is the weight connecting the i th hidden node to the output node.

The output function of ELM for generalized SLFNs can be identified by

$$f_N(\mathbf{x}) = \sum_{i=1}^N \beta_i G(\mathbf{a}_i, b_i, \mathbf{x}) = \beta \times h(\mathbf{x}) \quad (2)$$

For the binary classification applications, the decision function of ELM becomes

$$f_N(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N \beta_i G(\mathbf{a}_i, b_i, \mathbf{x}) \right) = \text{sign}(\beta \times h(\mathbf{x})) \quad (3)$$

Equation 2 can be written in another form as

$$H\beta = T \quad (4)$$

H and T are respectively hidden layer matrix and output matrix. Hidden layer matrix can be described as

$$H(\tilde{a}, \tilde{b}, \tilde{x}) = \begin{bmatrix} G(a_1, b_1, x_1) & \cdots & G(a_L, b_L, x_1) \\ \vdots & \ddots & \vdots \\ G(a_1, b_1, x_N) & \cdots & G(a_L, b_L, x_N) \end{bmatrix}_{N \times L} \quad (5)$$

where $\tilde{a} = a_1, \dots, a_L$, $\tilde{b} = b_1, \dots, b_L$, $\tilde{x} = x_1, \dots, x_N$. Output matrix can be described as

$$T = [t_1 \dots t_N]^T \quad (6)$$

The hidden nodes of SLFNs can be randomly generated. They can be independent of the training data.

3.2 AdaBoost

The AdaBoost [27] is a supervised learning algorithm designed to solve classification problems [28]. The algorithm takes as input a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ where the input sample $\mathbf{x}_i \in \mathbb{R}^p$, and the output value, y_i , in a finite space $y \in 1, \dots, K$. AdaBoost algorithm assumes, like ELM, a set of training data sampled independently and identically distributed (i.i.d.) from some unknown distribution \mathcal{X} .

Given a space of feature vectors X and two possible class labels, $y \in \{-1, +1\}$, AdaBoost goal is to learn a strong classifier $H(\mathbf{x})$ as a weighted ensemble of weak classifiers $h_t(\mathbf{x})$ predicting the label of any instance $\mathbf{x} \in X$ [29].

$$H(\mathbf{x}) = \text{sign}(f(\mathbf{x})) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right) \quad (7)$$

Pseudocode for AdaBoost is given in Alg. 1.

Algorithm 1 The AdaBoost algorithm.

```

1: Inputs:
    $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}, \mathbf{x}_i \in \mathbb{R}^p, y_i \in \{-1, +1\}\}_{i=1}^m$ 
2: Initialize  $\mathcal{D}(i) = \frac{1}{m}$  for all  $i$ 
3: while  $t < T$  do
4:   Train WeakLearner using distribution  $\mathcal{D}_t$ 
5:   get back a weak hypothesis  $h_t : X \rightarrow \{1, 2, \dots, K\}$ 
6:   calculate the error of  $h_t : \epsilon_t = Pr_{i \sim \mathcal{D}_t}[h_t(\mathbf{x}_i) \neq y_i]$ 
7:   Sets  $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$ 
8:   update distribution  $\mathcal{D}_{t+1} = \frac{\mathcal{D}_t}{Z_t} \times \begin{cases} e^{-\alpha_t}, & \text{if } h_t = y_i \\ e^{\alpha_t}, & \text{if } h_t \neq y_i \end{cases}$ 
9:   equivalently  $\mathcal{D}_{t+1} = \frac{\mathcal{D}_t \times \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$  where  $Z_t$  is a normalization
      constant.
10: end while
11: Outputs:
    final hypothesis  $h^* = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x})\right)$  .

```

3.3 MapReduce

MapReduce is a new programming model to run parallel applications for large scale data sets processing to support data-intensive applications. It is derived from the map and reduce function combination from functional programming. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. The MapReduce was originally developed by Google and built on principles in parallel manner [30]. The MapReduce framework first takes the input, divides it into smaller data chunks, and distributes them to worker nodes. MapReduce is divided into three major phases called map, reduce and a separated internal shuffle phase. The MapReduce framework automatically executes all those functions in a parallel manner over any number of processors/servers [31].

Pseudo code of MapReduce framework is shown in Eq. 8.

$$\begin{aligned}
 \text{map}(\text{key}_1, \text{value}_1) &\rightarrow \text{list}(\text{key}_2, \text{value}_2) \\
 \text{reduce}(\text{key}_2, \text{list}(\text{value}_2)) &\rightarrow \text{list}(\text{key}_3, \text{value}_3)
 \end{aligned} \tag{8}$$

Mapreduce programming technique is widely used on different scientific fields, i.e. cyber-security [32, 33], high energy physics [34], biology [35].

4 Proposed Approach

In this section we provide the details of the MapReduce based distributed AdaBoosted ELM algorithm. The basic idea of AdaBoost-ELM based on MapRe-

Table 2: Commonly used variables and notations.

Variables/Notation	Description
M	Data chunk split size
h	A single classifier function
X_m	Data chunk m of <i>input</i> values of \mathcal{D}
Y_m	Data chunk m of <i>output</i> values of \mathcal{D}
ϵ	Error rate
# Chunk	Number of data chunk
T	AdaBoost T size
# H. Nodes	Number of hidden nodes used in ELM
Acc	Accuracy of classifier hypothesis
k	Number of classes

duce technique is introduced in Section 4.1. The MapReduce implementation of AdaBoosted ELM is described in Section 4.3.

4.1 Basic Idea

Our main task is to parallel and distributed execute the computation of AdaBoosted ELM classification method. AdaBoosted ELM’s basic idea is to calculate ensemble of classifier functions over partitioned data (X_m, Y_m) in parallel manner. In Table 2, a summary of commonly used variables and notations to assess the classifier model performance of the AdaBoosted ELM method is given for convenience.

4.2 Analysis of the proposed algorithm

Barlett showed that the size of the weights is more important than the size of the neural network [36]. Kragh et al. also showed that ensemble methods of neural networks get better accuracy performance over unseen examples [37]. The main motivation of the this work is the idea that small size ELM ensembles can obtain more accurate classifier model that are comparable to individual classifiers.

In the proposed model, at every data chunk, there is a set of classifier functions that acts as a single classification model. The single model at every data chunk m is defined as follows:

$$f^{(m)}(\mathbf{x}) = \arg \max_k \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \quad (9)$$

The selected ensemble ELM classifier models from the reduce phase of MapReduce algorithm are combined into one single classification model.

$$\hat{h}(\mathbf{x}) = \arg \max_k \sum_{i=1}^m f^{(m)}(\mathbf{x}) \quad (10)$$

4.3 Implementation of the Model

The pseudocodes of MapReduce-based AdaBoost ELM are shown in Algorithm 2 and Algorithm 3. The *Map* procedure of our training model is implemented based on random assignment of each row of the training data set with split size of data, M , in line 2 of Algorithm 2. The input, \mathbf{x} , is a row of training data set \mathcal{D} . *Map* procedure partition the input matrix by row, producing $\langle randomSplitId, \mathbf{x} \rangle$ key-value pairs. *randomSplitId* is the identifier of the data chunk and is transferred as the input key to *Reduce* phase. The

Algorithm 2 AdaBoostedELM::Map

1: Inputs:

Training record $(\mathbf{x}, y) \in \mathcal{D}$, Data set split size M

2: $k \leftarrow rand(0, M)$

3: *Output*($k, (\mathbf{x}, y)$)

pseudo code of *Reduce* phase is shown in Algorithm 3. *Reduce* procedure is implemented based on the for-loop of lines 3 - 8 of Algorithm 3. The output ELM classifier of sub data set $(\mathbf{X}_k, \mathbf{y}_k)$ is calculated using AdaBoost constantly block by block, so every reduce task completes training phase and outputs an AdaBoosted set of classifier functions. The *mapper's* input k is the *randomSplitId* to create the data chunk and created in the *Map* phase of our training model.

Algorithm 3 AdaBoostedELM::Reduce

1: Inputs:

Key k , Value Set V , AdaBoost Iteration Size T

2: Split V into input space \mathbf{X}_n and out labels \mathbf{y}_n with $(\mathbf{X}_n, \mathbf{y}_n) \leftarrow V$

3: **for** $t = 1..T$ **do**

4: Train sub data set with ELM: $h_t \leftarrow ELM(\mathbf{X}, y)$

5: $\mathbf{y}_{pred}, \epsilon_t \leftarrow h_t(X)$

6: $\alpha_t \leftarrow \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$

7: $\mathcal{D}_{t+1} = \frac{\mathcal{D}_t \times \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$

8: **end for**

9: Outputs:

Final hypothesis for the reduce function $m :$

$$h_m \leftarrow \arg \max_k \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$$

5 Experiments

In this section, we perform experiments on real-world data sets from the public available data set repositories. Public data sets are used to evaluate the proposed learning method. Then, classification models of each data set are compared for accuracy results with the single instance of learning algorithm performance.

In Section 5.1 we explain the data sets and parameters that are used in experiments. The conventional ELM is applied all data sets and we find the accuracy performance over number of hidden nodes in Section 5.3. In Section 5.2, we show the empirical results of proposed distributed adaboost ELM training algorithm.

5.1 Experimental setup

In this section we apply our approach to five different data sets to verify its effectivity and efficiency. To demonstrate the effectiveness and performance of the proposed model, we apply it on various classification data sets from public data set repositories. To obtain an optimal value of Mapper size, m , we range it in the range from 20 to 100.

5.1.1 Commonly Used Classification Data Sets

We experiment on five public data sets which are summarized in Table 3, including Pendigit, Letter, Statlog, Page-blocks and Waveform. They are all multiclass data sets. All experiments are repeated 5 times and the results are averaged. All data sets are publicly available in svmlight format on the LIBSVM web site [38].

Pendigit data set is a collection of pen-based recognition of handwritten digits [39]. The data set contains 250 samples from 44 people. The first 7494 instances written by 30 people are used for the training data set, and the digits written by other 14 people are used for the independent testing purpose.

Skin data set is a collection of skin segmentation constructed over R, G, B color space [40]. The data set contains face images of different age groups (young, middle, old), genders and racial groups (White, Black, Asian). The data set contains 245057 instances; out of which 50859 is the skin labeled instances and 194198 is non-skin instances.

Statlog / Shuttle data set is a collection of space shuttle created by NASA [41]. The data set contains 43500 training instances and 14500 testing instances. 80% of the data belongs to class 1.

Page Blocks data set is a collection of page layout of a document that has been detected by a segmentation process [42]. The data set contains 4500 training instances and 973 testing instances.

Waveform data set is a collection of Breiman’s waveform domains of CART book’s [43]. The data set contains 4400 training instances and 600 testing instances.

Table 3: Description of the testing data sets used in the experiments.

Data set	#Train	#Test	#Classes	#Attributes
Pendigit	7494	3498	10	16
Skin	220543	24507	2	3
Statlog / Shuttle	43500	14500	7	9
Page-blocks	4500	973	5	10
Waveform	4400	600	3	21

5.1.2 Large Scale Classification Data Sets

We experiment on three public large scale data sets which are summarized in Table 4, including ”*Record Linkage Comparison Patterns (Donation)*”, ”*SUSY*” and ”*HIGGS*”. All experiments are repeated 5 times and the results are averaged.

Donation represent individual data, including first and family name, sex, date of birth and postal code, which were collected through iterative insertions in the course of several years. The comparison patterns in this data set are based on a sample of 100.000 records dating from 2005 to 2008 [44]. The data set contains 5,749,132 training instances and 1,000,000 testing instances. The data set is available on UCI web site [45].

SUSY is a classification data set that distinguish between a signal process which produces supersymmetric particles and a background process which does not [46]. The first 8 features are kinematic properties measured by the particle detectors in the accelerator. The last ten features are functions of the first 8 features. The data set contains 5,000,000 training instances and 50,000 testing instances. The data set is available on UCI web site [47].

HIGSS is a classification problem to distinguish between a signal process which produces Higgs bosons and a background process which does not [46]. The first 21 features (columns 2-22) are kinematic properties measured by the particle detectors in the accelerator. The last seven features are functions of the first 21 features. The data set contains 11,000,000 training instances and 500,000 testing instances. The data set is available on UCI web site [48].

5.2 Evaluation

Since the data sets that are used in our experiments are highly imbalanced, traditional accuracy based performance evaluation is not enough to find out

Table 4: Description of the testing large scale data sets used in the experiments.

Data set	#Train	#Test	#Classes	#Attributes
Donation	5,749,132	1,000,000	2	12
SUSY	5,000,000	50,000	2	18
HIGSS	11,000,000	1,000,000	2	28

an optimal classifier. We used four different metrics, the overall prediction accuracy, average recall, average precision [49] and F -score, to evaluate the classification accuracy which are common measurement metrics in information retrieval [50,51].

Precision is defined as the fraction of retrieved samples that are relevant. Precision is shown in Eq. 11.

$$Precision = \frac{Correct}{Correct + False} \quad (11)$$

Recall is defined as the fraction of relevant samples that is retrieved. Recall is shown in Eq. 12.

$$Precision = \frac{Correct}{Correct + Missed} \quad (12)$$

The proposed evaluation model calculates the precision and recall for each class from prediction scores then finds their mean. Average precision and recall is shown in Eq. 13 and Eq. 14.

$$Precision_{avg} = \frac{1}{n_{classes}} \sum_{i=0}^{n_{classes}-1} Prec_i \quad (13)$$

$$Recall_{avg} = \frac{1}{n_{classes}} \sum_{i=0}^{n_{classes}-1} Recall_i \quad (14)$$

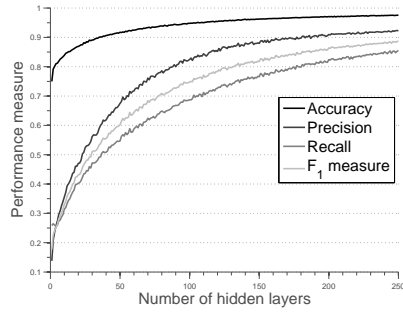
F -measure is defined as the harmonic mean of precision and recall. The

$$F_1 = 2 \times \frac{Prec_{avg} \times Recall_{avg}}{Prec_{avg} + Recall_{avg}} \quad (15)$$

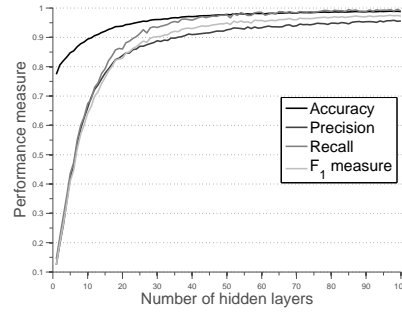
5.3 Data set results with conventional ELM

Figure 1 shows that the accuracy performance of ELM for experimental data sets becomes steady-state after a threshold value of N . The testing classification performance is measured through accuracy, precision, recall and F_1 measure. N varies from 150 to 500.

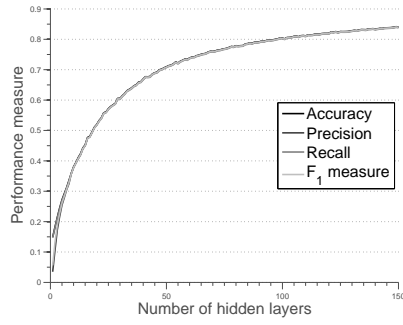
Table 5 shows the best performance of the conventional ELM method of each data set.



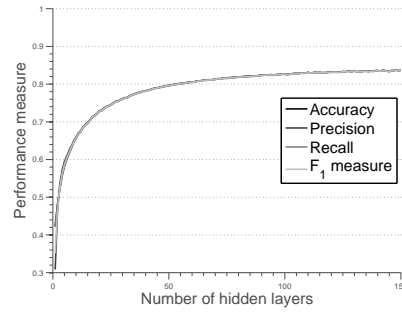
(a) Statlog data set.



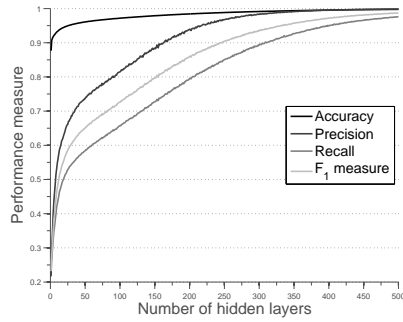
(b) Skin data set.



(c) Pen digit data set.



(d) Waveform data set.



(e) Page blocks data set.

Fig. 1: Number of hidden nodes in ELM versus classifier precision.

The conventional ELM training algorithm can be applied only in Section 5.1.1. The large scale data sets in Section 5.1.2 are not feasible to train on a single computer.

Table 5: Data set results with conventional ELM.

Data set	F_1	Recall	Precision	Accuracy
Pendigit	0.8404	0.8393	0.8416	0.8407
Skin	0.9754	0.9956	0.9583	0.9894
Statlog	0.8871	0.8556	0.9237	0.9757
Page-blocks	0.9873	0.9764	0.9988	0.9977
Waveform	0.8372	0.8368	0.8375	0.8376

5.4 Testing Accuracy Analysis

Because of two different data set type ("commonly used", "large scale") are used, the results are divided into two different sections. In Section 5.4.1, the figures and the plots show the implementation results of commonly used classification data sets. Section 5.4.2 shows the large scale data sets results.

5.4.1 Commonly Used Classification Data Sets

The results of accuracy and performance tests with real data are shown in Table 6 and Figure 2 - Figure 6. According to the these results, AdaBoost T size and Mapper size have more impact on the accuracy of ensemble ELM classifier than number of hidden nodes in ELM network.

Accuracy of classification models are visualized by heatmap color coding according to

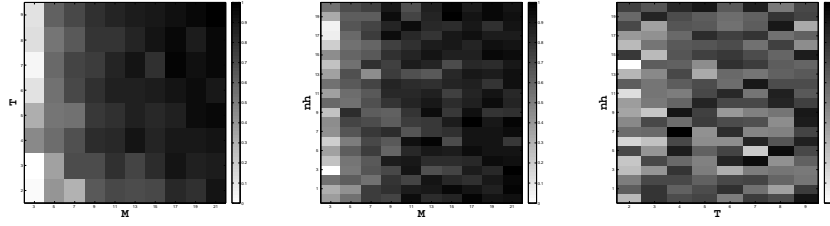
- Map size (M) - AdaBoost size (T)
- Map size (M) - Number of hidden nodes (nh)
- AdaBoost size (T) - Number of hidden nodes (nh)

Figure 2 - Figure 6 are used to plot the quantitative differences in accuracy score of each data set. Heatmaps are two dimensional graphical representations of data with a pre-defined colormap to display values of a matrix [52]. Heatmaps can be used to understand what parameters affect the accuracy of the classification model. The figures are used to comparatively illustrate accuracy levels across a number of different parameters including Map size, AdaBoost size and the number of hidden nodes in ELM algorithm obtained from the proposed learning method.

According to Table 7, classification performance results of the proposed method have almost the same values with the conventional ELM method.

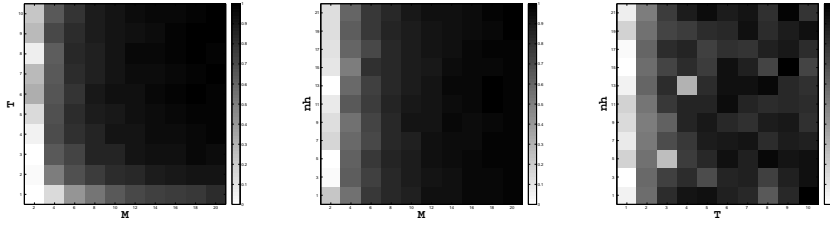
5.4.2 Large Scale Classification Data Sets

Figure 7 shows the speed up on mapper size over proposed method on large scale data sets. To asses the effectiveness of the learning algorithm, the time is measured with varying mapper size. Because of high dimensionality, the



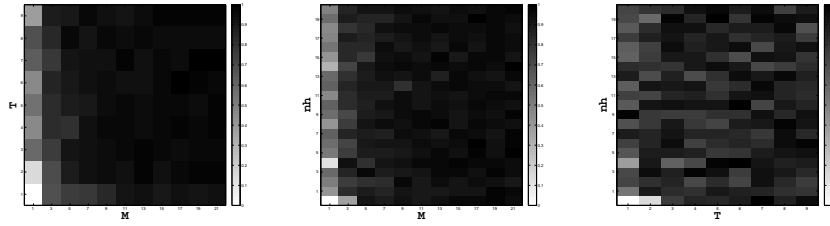
(a) Split size and adaboost T size (b) Split size and number of nh . (c) Adaboost T size and number of nh .

Fig. 2: Statlog data set heatmap.



(a) Split size and adaboost T size (b) Split size and number of nh . (c) Adaboost T size and number of nh .

Fig. 3: Pendigit data set heatmap.

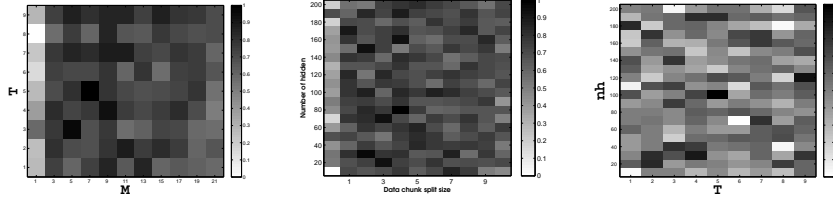


(a) Split size and adaboost T size (b) Split size and number of nh . (c) Adaboost T size and number of nh .

Fig. 4: Skin data set heatmap.

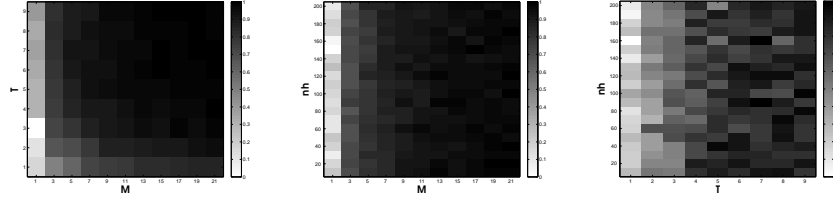
Table 6: Best performance results of data sets

Data set	# C.	T	# H.N.	Acc	Prec.	Recall	F_1
Pendigit	20	10	21	0,8256	0,8369	0,8234	0,8301
Skin	21	5	21	0,9892	0,9773	0,9913	0,9842
Statlog	11	2	21	0,9103	0,7486	0,5069	0,6045
Page Blocks	1	1	340	0,9404	0,9027	0,5756	0,7030
Waveform	19	6	40	0,862	0,8680	0,8605	0,8642



(a) Split size and adaboost T size (b) Split size and number of nh . (c) Adaboost T size and number of nh .

Fig. 5: Page blocks data set heatmap.



(a) Split size and adaboost T size (b) Split size and number of nh . (c) Adaboost T size and number of nh .

Fig. 6: Waveform data set heatmap.

Table 7: Performance comparison of ELM and proposed model.

Data set	Method	F_1	Recall	Precision	Accuracy
Pendigit	Conventional	0.8404	0.8393	0.8416	0.8407
	Proposed	0.8301	0.8234	0.8369	0.8256
Skin	Conventional	0.9754	0.9956	0.9583	0.9894
	Proposed	0.9842	0.9913	0.9773	0.9892
Statlog	Conventional	0.8871	0.8556	0.9237	0.9757
	Proposed	0.6045	0.5069	0.7486	0.9103
Page-blocks	Conventional	0.9873	0.9764	0.9988	0.9977
	Proposed	0.7030	0.5756	0.9027	0.9404
Waveform	Conventional	0.8372	0.8368	0.8375	0.8376
	Proposed	0.8642	0.8605	0.8680	0.8620

data sets cannot be trained on a single computer. Then, the standart speed up percentage is modified such that:

$$S_p = \frac{t_{\arg \min m \in M}}{t_p} \quad (16)$$

where $t_{\arg \min m \in M}$ is the total time on minimum mapper that can be achieved to build a classifier model.

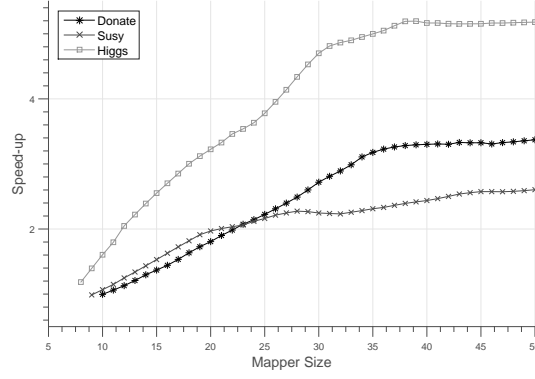


Fig. 7: Stability analysis of ensemble ELM classifiers with Mapper size.

As can be seen from the figure, the data sets achieves performance improvement in learning time of the algorithm. By examining the trends observed as the number of mappers increases, one can see that *non-linear* speed up is achieved.

5.5 Stability Analysis

Standard deviation of testing accuracy of the method is shown in Figure 8a and Figure 8b. We analyzed the stability of ensemble ELM classifier with two aspects, Mapper size and AdaBoost T size. Mapper size is the most important variable for the model stability according to the Figure 8a. From Figure 8a and Figure 8b, we can find that standard deviation of testing accuracy decreases enormously with the increasing of Mapper function size. Through this analysis, one can argue that a model with high Mapper function size do has higher stability than low Mapper function size.

6 Conclusion and Future Works

In this paper, a parallel AdaBoost extreme learning machine algorithm implementation has been proposed for massive data learning. By creating the overall data set into data chunks, MapReduce based learning algorithm reduces the training time of ELM classification. To overcome the accuracy performance decreasing, distributed ELM is enhanced with AdaBoost method. The experimental results show that AdaBoosted ELM not only reduce the training time of large-scale data sets, but also evaluation metrics of accuracy performance compared with the conventional ELM.

The proposed AdaBoost based ELM has three different trade-off parameters which are (i) data chunk split size, M , (ii) maximum number of iterations,

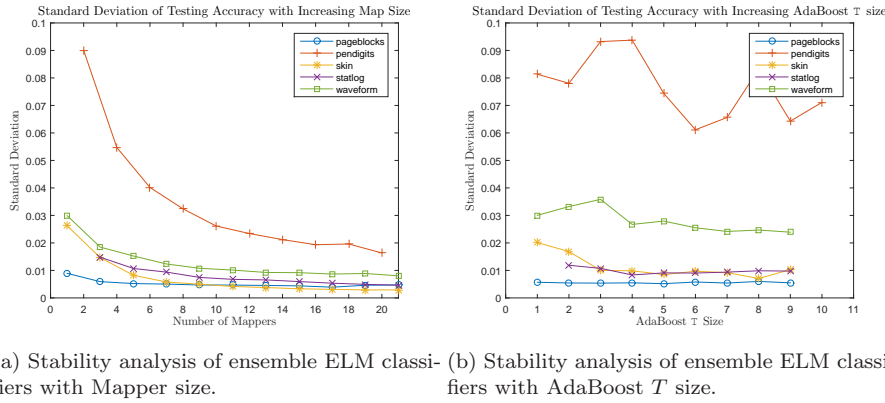


Fig. 8: Stability analysis

T , in AdaBoost Algorithm and lastly (iii) number of hidden layer nodes nh in ELM algorithm. The empirical results in heatmap figures show that parameters M and T are more dominant than parameter nh for the classification accuracy of the hypothesis.

The algorithm is designed to deal with large scale data set ELM training problems. Another objective is to achieve the model's classification performance with same or close to the conventional ELM method. Classification performance results are shown in Section 5.3. The empirical results show us that classification performance results of the proposed method have almost the same values with the conventional ELM method.

References

1. G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, "Extreme learning machine: Theory and applications," *Neurocomputing*, vol. 70, no. 13, pp. 489 – 501, 2006. Neural Networks Selected Papers from the 7th Brazilian Symposium on Neural Networks (SBRN '04) 7th Brazilian Symposium on Neural Networks.
2. X.-g. Zhao, G. Wang, X. Bi, P. Gong, and Y. Zhao, "Xml document classification based on elm," *Neurocomputing*, vol. 74, no. 16, pp. 2444–2451, 2011.
3. G. Wang, Y. Zhao, and D. Wang, "A protein secondary structure prediction framework based on the extreme learning machine," *Neurocomputing*, vol. 72, no. 13, pp. 262 – 268, 2008. Machine Learning for Signal Processing (MLSP 2006) / Life System Modelling, Simulation, and Bio-inspired Computing (LSMS 2007).
4. W. Zong and G.-B. Huang, "Face recognition based on extreme learning machine," *Neurocomputing*, vol. 74, no. 16, pp. 2541 – 2551, 2011. Advances in Extreme Learning Machine: Theory and Applications Biological Inspired Systems. Computational and Ambient Intelligence Selected papers of the 10th International Work-Conference on Artificial Neural Networks (IWANN2009).
5. Y. Lan, Z. Hu, Y. C. Soh, and G.-B. Huang, "An extreme learning machine approach for speaker recognition," *Neural Computing and Applications*, vol. 22, no. 3-4, pp. 417–425, 2013.
6. Y. Lu, V. Roychowdhury, and L. Vandenberghe, "Distributed parallel support vector machines in strongly connected networks," *Neural Networks, IEEE Transactions on*, vol. 19, pp. 1167–1178, July 2008.

7. Z. Sun and G. Fox, "Study on parallel svm based on mapreduce," in *International Conference on Parallel and Distributed Processing Techniques and Applications*, pp. 16–19, Citeseer, 2012.
8. F. Catak and M. Balaban, "Cloudsvm: Training an svm classifier in cloud computing systems," in *Pervasive Computing and the Networked World* (Q. Zu, B. Hu, and A. Eli, eds.), vol. 7719 of *Lecture Notes in Computer Science*, pp. 57–68, Springer Berlin Heidelberg, 2013.
9. B. Panda, J. S. Herbach, S. Basu, and R. J. Bayardo, "Planet: Massively parallel learning of tree ensembles with mapreduce," *Proc. VLDB Endow.*, vol. 2, pp. 1426–1437, Aug. 2009.
10. C. Zhang, F. Li, and J. Jests, "Efficient parallel knn joins for large data in mapreduce," in *Proceedings of the 15th International Conference on Extending Database Technology, EDBT '12*, (New York, NY, USA), pp. 38–49, ACM, 2012.
11. T. Sun, C. Shu, F. Li, H. Yu, L. Ma, and Y. Fang, "An efficient hierarchical clustering method for large datasets with map-reduce," in *Parallel and Distributed Computing, Applications and Technologies, 2009 International Conference on*, pp. 494–499, Dec 2009.
12. Y. He, H. Tan, W. Luo, H. Mao, D. Ma, S. Feng, and J. Fan, "Mr-dbscan: An efficient parallel density-based clustering algorithm using mapreduce," in *Parallel and Distributed Systems (ICPADS), 2011 IEEE 17th International Conference on*, pp. 473–480, Dec 2011.
13. W. Zhao, H. Ma, and Q. He, "Parallel k-means clustering based on mapreduce," in *Cloud Computing* (M. Jaatun, G. Zhao, and C. Rong, eds.), vol. 5931 of *Lecture Notes in Computer Science*, pp. 674–679, Springer Berlin Heidelberg, 2009.
14. J. Xin, Z. Wang, C. Chen, L. Ding, G. Wang, and Y. Zhao, "Elm: distributed extreme learning machine with mapreduce," *World Wide Web*, vol. 17, no. 5, pp. 1189–1204, 2014.
15. N.-Y. Liang, G.-B. Huang, P. Saratchandran, and N. Sundararajan, "A fast and accurate online sequential learning algorithm for feedforward networks," *Neural Networks, IEEE Transactions on*, vol. 17, pp. 1411–1423, Nov 2006.
16. Y. Sun, Y. Yuan, and G. Wang, "An os-elm based distributed ensemble classification framework in {P2P} networks," *Neurocomputing*, vol. 74, no. 16, pp. 2438 – 2443, 2011. Advances in Extreme Learning Machine: Theory and Applications Biological Inspired Systems. Computational and Ambient Intelligence Selected papers of the 10th International Work-Conference on Artificial Neural Networks (IWANN2009).
17. B. Wang, S. Huang, J. Qiu, Y. Liu, and G. Wang, "Parallel online sequential extreme learning machine based on mapreduce," *Neurocomputing*, vol. 149, Part A, no. 0, pp. 224 – 232, 2015. Advances in neural networks Selected papers from the Tenth International Symposium on Neural Networks (ISNN 2013) Advances in Extreme Learning Machines Selected articles from the International Symposium on Extreme Learning Machines (ELM 2013).
18. X. Bi, X. Zhao, G. Wang, P. Zhang, and C. Wang, "Distributed extreme learning machine with kernels based on mapreduce," *Neurocomputing*, vol. 149, Part A, no. 0, pp. 456 – 463, 2015. Advances in neural networks Selected papers from the Tenth International Symposium on Neural Networks (ISNN 2013) Advances in Extreme Learning Machines Selected articles from the International Symposium on Extreme Learning Machines (ELM 2013).
19. J. Chen, G. Zheng, and H. Chen, "Elm-mapreduce: Mapreduce accelerated extreme learning machine for big spatial data analysis," in *Control and Automation (ICCA), 2013 10th IEEE International Conference on*, pp. 400–405, June 2013.
20. L. I. Kuncheva and C. J. Whitaker, "Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy," *Machine learning*, vol. 51, no. 2, pp. 181–207, 2003.
21. G. bin Huang, Q. yu Zhu, and C. kheong Siew, "Extreme learning machine: A new learning scheme of feedforward neural networks," in *IN PROC. INT. JOINT CONF. NEURAL NETW*, pp. 985–990, 2006.
22. G.-B. Huang, L. Chen, and C.-K. Siew, "Universal approximation using incremental constructive feedforward networks with random hidden nodes," *Neural Networks, IEEE Transactions on*, vol. 17, pp. 879–892, July 2006.

23. G.-B. Huang and L. Chen, "Convex incremental extreme learning machine," *Neurocomputing*, vol. 70, no. 1618, pp. 3056 – 3062, 2007. Neural Network Applications in Electrical Engineering Selected papers from the 3rd International Work-Conference on Artificial Neural Networks (IWANN 2005) 3rd International Work-Conference on Artificial Neural Networks (IWANN 2005).
24. G.-B. Huang and L. Chen, "Enhanced random search based incremental extreme learning machine," *Neurocomputing*, vol. 71, no. 1618, pp. 3460 – 3468, 2008. Advances in Neural Information Processing (ICONIP 2006) / Brazilian Symposium on Neural Networks (SBRN 2006).
25. J. Tang, C. Deng, G.-B. Huang, and B. Zhao, "Compressed-domain ship detection on spaceborne optical image using deep neural network and extreme learning machine," *Geoscience and Remote Sensing, IEEE Transactions on*, vol. 53, pp. 1174–1185, March 2015.
26. G.-B. Huang, M.-B. Li, L. Chen, and C.-K. Siew, "Incremental extreme learning machine with fully complex hidden nodes," *Neurocomputing*, vol. 71, no. 46, pp. 576 – 583, 2008. Neural Networks: Algorithms and Applications 4th International Symposium on Neural Networks 50 Years of Artificial Intelligence: a Neuronal Approach Campus Multidisciplinary in Perception and Intelligence.
27. Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," in *Computational learning theory*, pp. 23–37, Springer, 1995.
28. Y. Freund, R. Schapire, and N. Abe, "A short introduction to boosting," *Journal-Japanese Society For Artificial Intelligence*, vol. 14, no. 771-780, p. 1612, 1999.
29. I. Landesa-Vzquez and J. L. Alba-Castro, "Double-base asymmetric adaboost," *Neurocomputing*, vol. 118, no. 0, pp. 101 – 114, 2013.
30. J. Dean and S. Ghemawat, "Mapreduce: Simplified data processing on large clusters," *Commun. ACM*, vol. 51, pp. 107–113, Jan. 2008.
31. M. C. Schatz, "CloudBurst: highly sensitive read mapping with MapReduce," *Bioinformatics (Oxford, England)*, vol. 25, pp. 1363–1369, June 2009.
32. J. Choi, C. Choi, B. Ko, and P. Kim, "A method of ddos attack detection using http packet pattern and rule engine in cloud computing environment," *Soft Computing*, vol. 18, no. 9, pp. 1697–1703, 2014.
33. M. Ogiela, A. Castiglione, and I. You, "Soft computing for security services in smart and ubiquitous environments," *Soft Computing*, vol. 18, no. 9, pp. 1655–1658, 2014.
34. W. Bhimji, T. Bristow, and A. Washbrook, "Hepdoop: High-energy physics analysis using hadoop," in *Journal of Physics: Conference Series*, vol. 513, p. 022004, IOP Publishing, 2014.
35. L. Xu, H. Kim, X. Wang, W. Shi, and T. Suh, "Privacy preserving large scale dna read-mapping in mapreduce framework using fpgas," in *Field Programmable Logic and Applications (FPL), 2014 24th International Conference on*, pp. 1–4, IEEE, 2014.
36. P. Bartlett, "The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network," *Information Theory, IEEE Transactions on*, vol. 44, pp. 525–536, Mar 1998.
37. A. Krogh and J. Vedelsby, "Neural network ensembles, cross validation, and active learning," in *Advances in Neural Information Processing Systems*, pp. 231–238, MIT Press, 1995.
38. LIBSVM, "Libsvm data: Classification, regression, and multi-label." <http://ntu.edu.tw/~csie/ntu.edu.tw/>
39. F. Alimoglu and E. Alpaydin, "Methods of combining multiple classifiers based on different representations for pen-based handwritten digit recognition," in *Proceedings of the Fifth Turkish Artificial Intelligence and Artificial Neural Networks Symposium (TAINN 96)*, 1996.
40. R. Bhatt, G. Sharma, A. Dhall, and S. Chaudhury, "Efficient skin region segmentation using low complexity fuzzy decision tree model," in *India Conference (INDICON), 2009 Annual IEEE*, pp. 1–4, Dec 2009.
41. C.-W. Hsu and C.-J. Lin, "A comparison of methods for multiclass support vector machines," *Trans. Neur. Netw.*, vol. 13, pp. 415–425, Mar. 2002.

42. D. Malerba, F. Esposito, and G. Semeraro, "A further comparison of simplification methods for decision-tree induction," in *In D. Fisher and H. Lenz (Eds.), Learning*, pp. 365–374, Springer-Verlag, 1996.
43. L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen, *Classification and regression trees*. CRC press, 1984.
44. I. Schmidtman, G. Hammer, M. Sariyar, A. Gerhold-Ay, and K. des öffentlichen Rechts, "Evaluation des krebregisters nrw–schwerpunkt record linkage," *Abschlußbericht vom*, vol. 11, 2009.
45. UCI, "Record linkage comparison patterns data set," 2011. Available online at <https://archive.ics.uci.edu/ml/datasets/Record+Linkage+Comparison+Patterns>.
46. P. Baldi, P. Sadowski, and D. Whiteson, "Searching for exotic particles in high-energy physics with deep learning," *Nature communications*, vol. 5, 2014.
47. UCI, "Susy data set," 2014. Available online at <https://archive.ics.uci.edu/ml/datasets/SUSY>.
48. UCI, "Higgs data set," 2014. Available online at <https://archive.ics.uci.edu/ml/datasets/HIGGS>.
49. A. Turpin and F. Scholer, "User performance versus precision measures for simple search tasks," in *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '06, (New York, NY, USA), pp. 11–18, ACM, 2006.
50. C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press, 2008.
51. J. Makhoul, F. Kubala, R. Schwartz, and R. Weischedel, "Performance measures for information extraction," in *In Proceedings of DARPA Broadcast News Workshop*, pp. 249–252, 1999.
52. B. Khomtchouk, D. Van Booven, and C. Wahlestedt, "Heatmapgenerator: high performance rnaseq and microarray visualization software suite to examine differential gene expression levels using an r and c++ hybrid computational pipeline," *Source Code for Biology and Medicine*, vol. 9, no. 1, 2014.